



1

# Exploring Novel Neuromorphic Computing Architectures with a Multi-Node FPGA System

**Pierpaolo Perticaroli** INFN Roma 1 – APE Lab

Workshop sul Calcolo nell'INFN 26-30 May 2025



### BRAINSTAIN - CSN5 PROJECT BRAIN Studies and Technologies for Artificial Intelligence and Neuroscience





### BRAINSTAIN - CSN5 PROJECT BRAIN Studies and Technologies for Artificial Intelligence and Neuroscience





# **Neuromorphic Computing - Why mimic the Brain?**



In-Memory Computing. Co-location of processing and storage

#### Traditional Von Neuman computing. The cost of data movement

Operation	Energy consumption
Addition of data	1x
Access data (onchip cache)	60x
Access data (offchip RAM)	3500x



Brain. Closely tied processing-storage, highly parallel



Sze et al, IEEE Custom Integrated Circuits Conference (2017)

- Sparse, distributed information encoding through spikes
- Asynchronous, event-based computation relying on spikes time
- Local learning rules
  - No costly back-propagation

### High energy efficiency, real-time processing, incremental learning





# **Spiking Neural Networks Applications**



#### from snnTorch **Rate Coding** Latency Coding **Delta Modulation Energy efficiency**, nput = 0 **Real-time learning** B(n=1, p=0) = 0Input Data Input Data input = 0.5B(n=1, p=0.5) = ?Ideal for edge Computing, on-sensor Onspikes processing, dataflow applications input = 1 B(n=1, p=1) = 1 Any input can be encoded as On spikes + off spikes Data is repeated spikes along the first Each feature corresponds to a Input features are used to dimension for single spike. The intensity of parameterize a binomial Particularly effective when rate coding, i.e., the feature determines how distribution, which is then time-first, fast the spike occurs. Options sampled from to determine before spike for linear or logarithmic firing input is already sparse, eventwhether or not a spike occurs encoding times are available. based or is acquired as a time Si Pixel Sensor sequence **SNN** Identify interesting applications in physics (particle physics sensors, anomaly detection, event Charge cameras, time series data ...)



### What we are working on



#### Multi-node FPGA system

 Modular, scalable and reconfigurable neuromorphic architecture

### Our target

Accelerator platform for Spiking NN

Edge and dataflow applications in physics experiments



### What we are working on



#### Multi-node FPGA system

- Modular, scalable and reconfigurable neuromorphic architecture
- Support for biologically realistic neuron dynamics
- Support for multi-compartment neuron models
- Deterministic, reproducible results

### Our target

Accelerator platform for Spiking NN

Edge and dataflow applications in physics experiments

 Neuroscience simulator for specific use cases

Explore application of novel bio-inspired AI



## What we are working on



#### Multi-node FPGA system

- Modular, scalable and reconfigurable neuromorphic architecture
- Support for biologically realistic neuron dynamics
- Support for multi-compartment neuron models



### Our target

Accelerator platform for Spiking NN

Edge and dataflow applications in physics experiments

 Neuroscience simulator for specific use cases

Explore application of novel bio-inspired AI

#### Focus on

Simplified two-compartment neuron with calcium dynamics capturing brainstate specific apical-amplification, isolation and -drive Pastorelli et Al. 2025



# **Architecture Design – Workflow and Tools**





- SystemC (C++ based event-driven) architecture simulator
  - Validate and test architecture functionality
  - Modelling multi-FPGAs communication delays
  - SystemC allows higher or lower levels of abstraction (TLM, RTL-like)

#### Catapult High Level Synthesis (HLS)

- Direct transition from high-level simulator to RTL
- Benefits: high level language, software emulation mode, early estimates on latency/throughput and FPGA-resources consumption → relatively easy validation, reprogrammability, and debug

#### Catapult HLS, over Vitis HLS

- Supports SystemC Synthesis
- Vitis HLS C++ Software Emulation is not intuitevily applicable to a multi-FPGA simulation model
- SystemC RTL-like flexibility is useful to model architectural details such as synchronization and feedback between hardware blocks



### **Architecture Design**



#### INFN Comm. IP



- Direct network of interconnected FPGAs
- Baseline topology assumption: 3D Torus
- Relying on internal INFN Communication IP for lowlatency inter-node and intra-node communication (< 1 us for up to 1 kB packets)
  - Node = FPGA
- NEURO-CORE
  - Core architectural component, containing a set of neurons
  - Exchanges spikes and barrier synchronization messages with other Neuro-Cores through the INFN Comm.IP
  - Developed with High Level Synthesis (HLS) tools
  - Switching to lower level VHDL if needed for performance of critical blocks



### **Architecture Design**













#### **Modular architecture**

- Neuro-Core contains Workers (Wj)
- Workers contain neurons (n)

#### Reconfigurable

- Optimal Bit-filed sizes and parameters (Nr of Workers, Nr of Neurons per Worker) will depend on target application
- Synaptic Memory can be implemented as offchip DRAM or local BRAM, according to size of target network connectome





#### **Modular architecture**

- Neuro-Core contains Workers (Wj)
- Workers contain neurons (n)

#### Reconfigurable

- Optimal Bit-filed sizes and parameters (Nr of Workers, Nr of Neurons per Worker) will depend on target application
- Synaptic Memory can be implemented as offchip DRAM or local BRAM, according to size of target network connectome





#### Spike input

Characterized by ID of pre-synaptic neuron who sent it Address Event Representation (AER)





### Dataflow





#### Memory Look-Up and distribution

- Memory Look-Up according to **ID**<sub>pre</sub>
  - Extract the weight and delay of synapses connecting the  $\mathsf{ID}_{\mathsf{pre}}$  neuron to the post-synaptic neurons on this core
- Distribute this data to the post-synaptic neurons
  - Multiple out lanes to increase throughput





#### Memory Look-Up and distribution

- Memory Look-Up according to **ID**<sub>pre</sub>
  - Extract the weight and delay of synapses connecting the  $ID_{pre}$  neuron to the post-synaptic neurons on this core
- Distribute this data to the post-synaptic neurons
  - Multiple out lanes to increase throughput
  - Load balancing Round Robin distribution to local Workers containing post-synaptic neurons



0

\_Synaptic

data



- When spike synaptic data (weight,delay) reaches here, the weight gets collocated in ring buffer according to delay, to be processed at a later time, mimicking synaptic delays
- Summed to previous spikes' weights already in that time slot, making for an aggregate weight







 When all pre-synaptic spikes have been received

#### Worker can start Neurons

#### **Dynamics**

Currently Leaky Integrate and Fire with delta shaped post synaptic currents (*iaf\_psc\_delta* from NEST simulator)

$$rac{dV_{
m m}}{dt} = -rac{V_{
m m}-E_{
m L}}{ au_{
m m}} + \dot{\Delta}_{
m syn} + rac{I_{
m syn}+I_{
m e}}{C_{
m m}}$$

$$\dot{\Delta}_{
m syn}(t) = \sum_j w_j \sum_k \delta(t-t_j^k-d_j)$$

HLS allows agile study of different neuron models









**Spike emission** 



# **Functional Validation**

 Comparing simulation results of our system with the enstablished NEST simulator on a common network testcase: same neuron model, connectivity, external input spike trains.

# nest::

C++ simulator, with a Python interface, implementing a large number of models of biological neurons and synapses

0 0 0 1 0.

- Bit-exact agreement expected for software-simulator but not for final hardware implementation for multiple reasons, e.g.:
   64-bit floating point arithmetic in NEST vs fixed-point
  - 64-bit floating point arithmetic in NEST vs fixed-point numeric representation in our system; currently 32-bit, but evaluating smaller precision.
  - Neuron dynamics implementation: differential equations solver methods suited for software vs hardware-optimized methods
  - MPI communication vs custom communication





# **Spike-train distance metrics**

Studying effects of reduced fixed-point representation, using the SystemC RTLlike bit-accurate simulator

#### Neuron mean Inter-spike Interval (ISI) difference FP – FXP32



### **FPGA Occupancy Single Neuro-core on a Versal VPK-180**

#### W = workers

n = neurons

VPK-180 Total PL mem 994 Mb

Single Neuro- Core	LUT	FF	DSP	BRAM	Synaptic Mem requires
16 n (4 W, 4 n x W)	13.5k <b>0.40%</b>	21.6k <b>0.32%</b>	0 <b>0.00%</b>	2 <b>0.04%</b>	16 Kb
64 n (4 W, 16 n x	21.1k	31.4k	0	4	197 Kb
W)	<b>0.63%</b>	<b>0.47%</b>	<b>0.00%</b>	<b>0.08%</b>	
512 n (4 W, 128 n x	75.5k	30.6k	0	16	4 Mb
W)	<b>2.25%</b>	<b>0.45%</b>	<b>0.00%</b>	<b>0.32%</b>	
2048 n (32 W, 64 n	764.0k	197.3k	2	96	84 Mb
x W)	<b>22.7%</b>	<b>2.93%</b>	<b>0.01%</b>	<b>1.94%</b>	
2048 n (16 W, 128 n	708.5k	106.9k	34	64	84 Mb
x W)	<b>21.1%</b>	<b>1.59%</b>	<b>0.24%</b>	<b>1.30%</b>	



- Currently validating complete hardware testbed, with a Neuro-Core + a Poisson input Core on a single FPGA Next steps
- Begin deployment of a target application, starting with a single FPGA
- Optimize fixed-point representation, architecture parameters and neuron model alongside the target application performance
- Validate multi-FPGA design